

Ten Years of Hunting for Similar Code for Fun and Profit

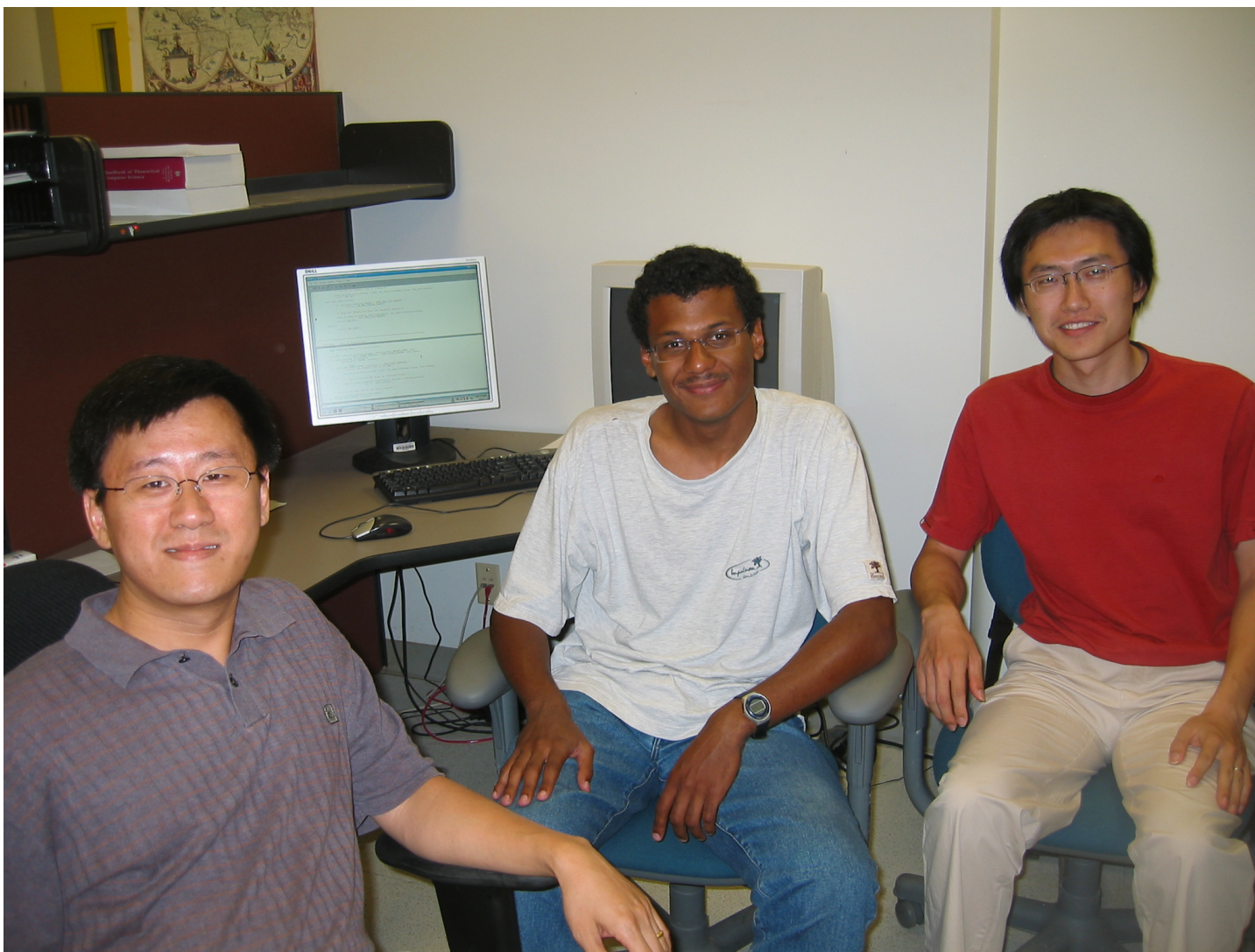
Stéphane Glondou
Lingxiao Jiang
Zhendong Su

Software is about similarity!

DECKARD: Scalable and Accurate Tree-based Detection of Code Clones*

Lingxiao Jiang Ghassan Misherghi Zhendong Su
University of California, Davis
`{lxjiang,ghassanm,su}@ucdavis.edu`

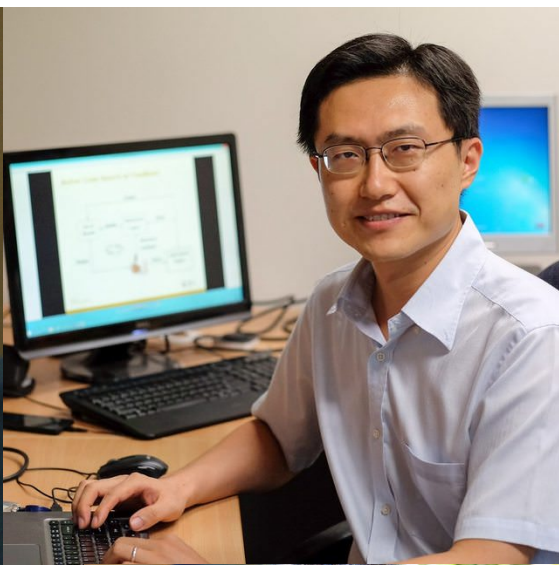
Stéphane Glondou
ENS de Cachan, France
`steph@glondou.net`



2005



2006



But, what is **impact**?



Concept

Insight

Technique

Tool

Software eats everything

Software is about similarity

Algorithms

Procedural
Abstraction

Algorithms

Procedural
Abstraction

Modularity

Algorithms

Procedural
Abstraction

Modularity

Design
Patterns

Algorithms

Procedural
Abstraction

Modularity

Design
Patterns

Object
Orientation

Algorithms

Procedural
Abstraction

Modularity

Design
Patterns

Object
Orientation

Algorithms

Refactoring

Procedural
Abstraction

Modularity

Design
Patterns

Object
Orientation

Algorithms

Refactoring

Software
Analytics

Procedural
Abstraction

Modularity

Design
Patterns

Object
Orientation

Algorithms

Refactoring

Software
Analytics

“Big Code”

How to **define** similarity?

Textual

Token

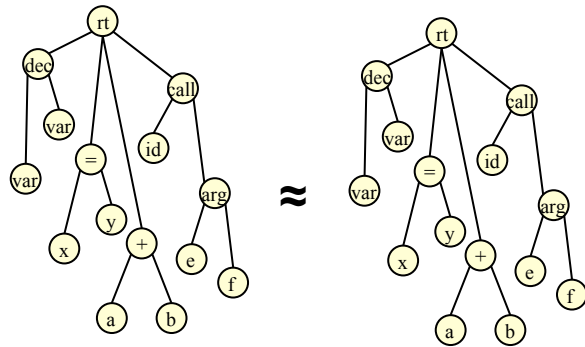
AST

PDG

Functional

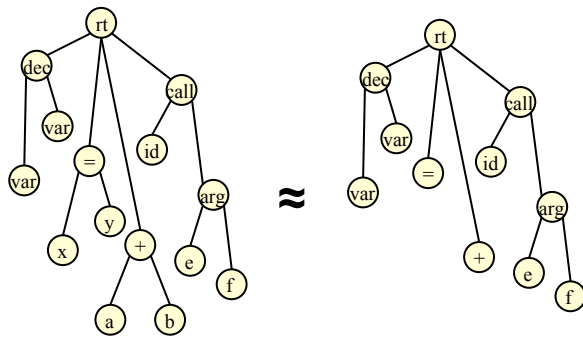
Deckard

Goal



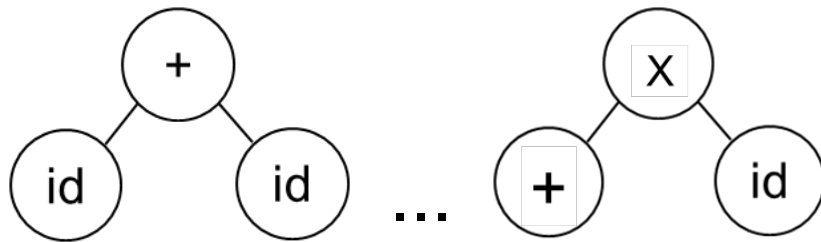
*Scalable tree
similarity*

Concept



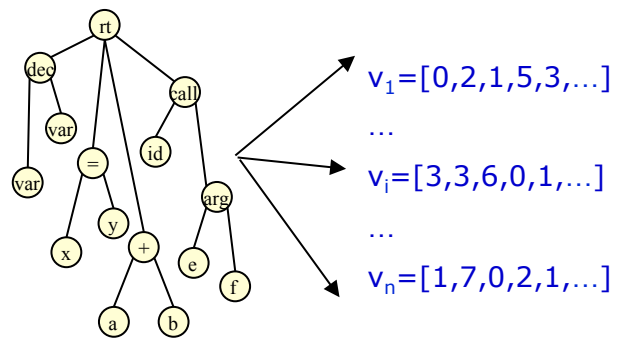
*Approximate
tree similarity*

Insight



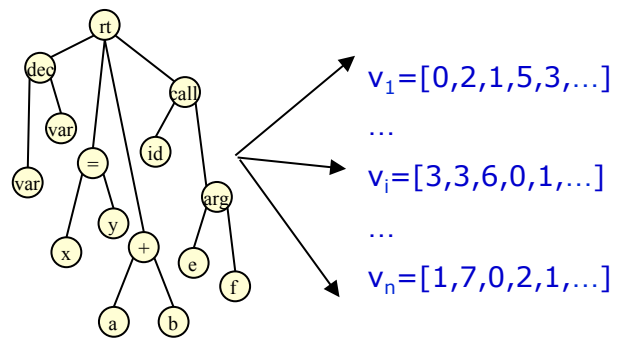
q-level subtrees

Technique



*Vector embedding
of trees*

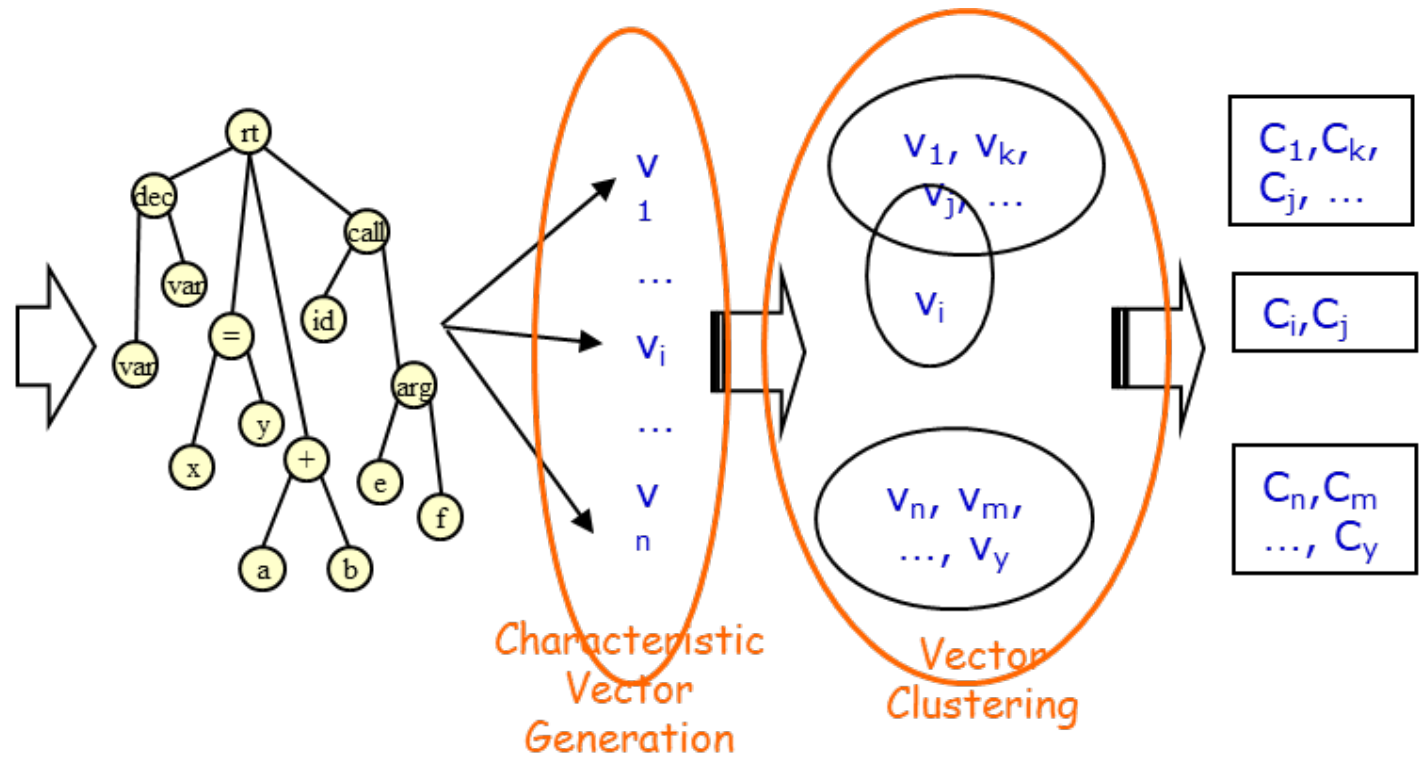
Technique



*Vector embedding
of trees*

*Locality sensitive
hashing*

Program

[illegible]

Tool



Deckard
(on GitHub)

Follow-up

ISSTA'09, July 19–23, 2009, Chicago, Illinois, USA

Detecting Code Clones in Binary Executables*

Andreas Sæbjørnsen
University of California, Davis
andsebj@ucdavis.edu

Jeremiah Willcock[†]
Indiana University
jewillco@osl.iu.edu

Thomas Panas
Lawrence Livermore National
Laboratory
panas2@llnl.gov

Daniel Quinlan
Lawrence Livermore National
Laboratory
dquinlan@llnl.gov

Zhendong Su
University of California, Davis
su@ucdavis.edu

ESEC-FSE'07: Proceedings of the 6th Joint Meeting of the European Software Engineering Conference and the ACM
SIGSOFT Symposium on the Foundations of Software Engineering: Dubrovnik, Croatia, September 3-7, 2007, Pages 55-64.
<http://doi.org/10.1145/1287624.1287634>

Context-Based Detection of Clone-Related Bugs*

Lingxiao Jiang

Zhendong Su

Edwin Chiu

Department of Computer Science
University of California, Davis
{lxjiang, su, eychiu}@ucdavis.edu

OOPSLA/SPLASH'10, October 17–21, 2010, Reno/Tahoe, Nevada, USA

Scalable and Systematic Detection of Buggy Inconsistencies in Source Code

Mark Gabel¹ Junfeng Yang^{2*} Yuan Yu³ Moises Goldszmidt³ Zhendong Su^{1†}

¹University of California, Davis
{mggabel, su}@ucdavis.edu

²Columbia University
junfeng@cs.columbia.edu

³Microsoft Research, Silicon Valley
{yuanbyu, moises}@microsoft.com

FSE '14, November 16–22, 2014, Hong Kong, China

Vector Abstraction and Concretization for Scalable Detection of Refactorings

Narcisa Andreea Milea
School of Computing
National University of
Singapore
mileanar@comp.nus.edu.sg

Lingxiao Jiang
School of Information Systems
Singapore Management
University
lxjiang@smu.edu.sg

Siau-Cheng Khoo
School of Computing
National University of
Singapore
khoosc@nus.edu.sg

ICSE'08, May 10–18, 2008, Leipzig, Germany

Scalable Detection of Semantic Clones*

Mark Gabel Lingxiao Jiang Zhendong Su

Department of Computer Science
University of California, Davis
{mggabel, lxjiang, su}@ucdavis.edu

ISSTA '09: Proceedings of the 2009 International Symposium on Software Testing and Analysis, Chicago, July 19-23, 2009, Pages 81-92.
<http://doi.org/10.1145/1572272.1572283>

Automatic Mining of Functionally Equivalent Code Fragments via Random Testing*

Lingxiao Jiang
University of California, Davis
jiangl@cs.ucdavis.edu

Zhendong Su
University of California, Davis
su@cs.ucdavis.edu

FSE-18, November 7–11, 2010, Santa Fe, New Mexico, USA

A Study of the Uniqueness of Source Code*

Mark Gabel Zhendong Su

Department of Computer Science
University of California at Davis
{mggabel,su}@ucdavis.edu

ICSE 2012, Zurich, Switzerland

On the Naturalness of Software

Abram Hindle, Earl T. Barr, Zhendong Su

Dept. of Computer Science

University of California at Davis

Davis, CA 95616 USA

{ajhindle,barr,su}@cs.ucdavis.edu

Mark Gabel

Dept. of Computer Science

The University of Texas at Dallas

Richardson, TX 75080 USA

mark.gabel@utdallas.edu

Premkumar Devanbu

Dept. of Computer Science

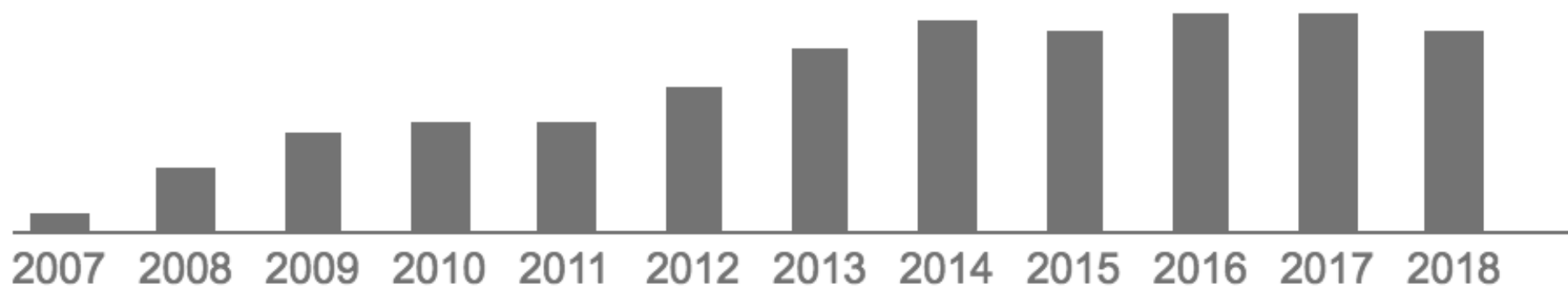
University of California at Davis

Davis, CA 95616 USA

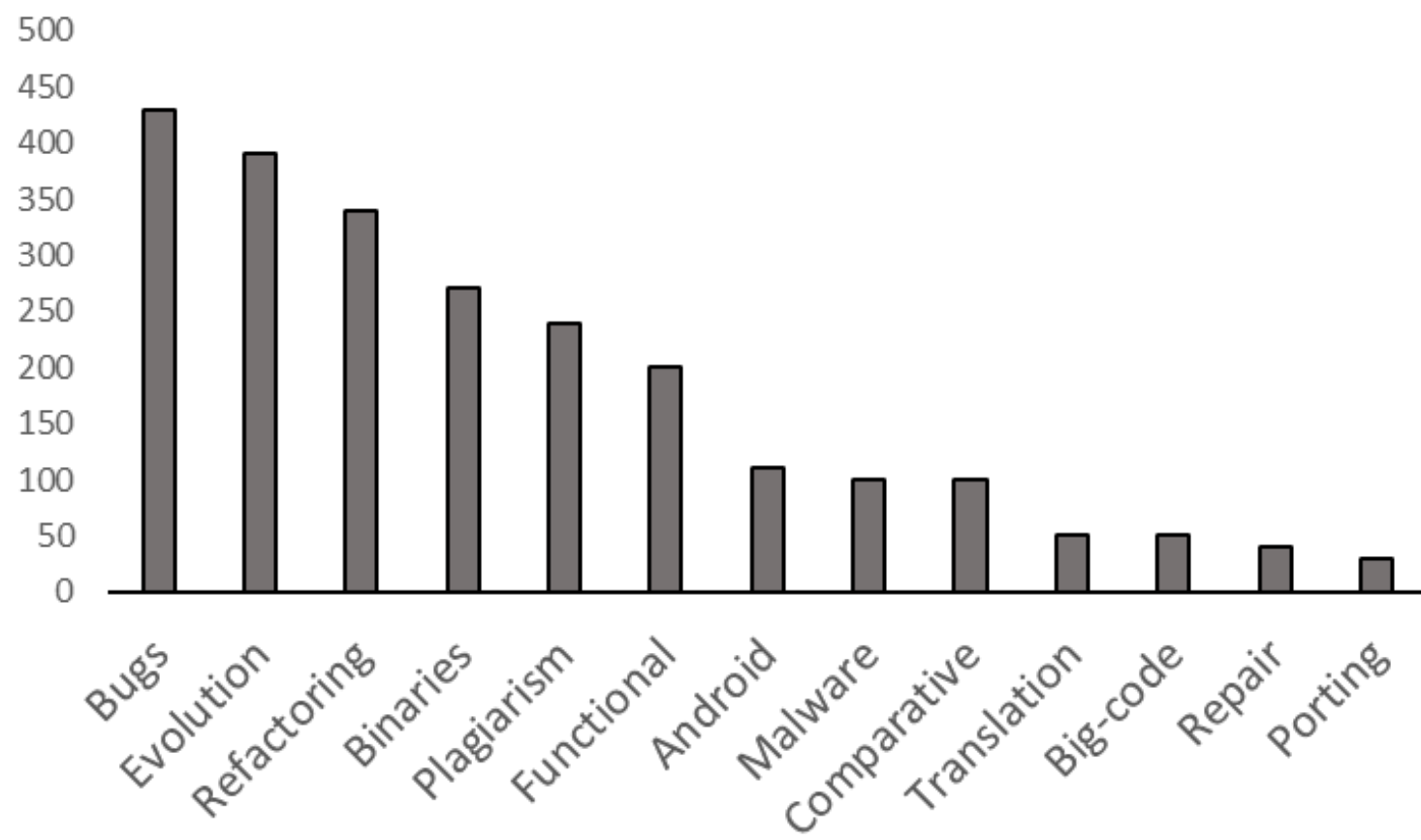
devanbu@cs.ucdavis.edu

Influence

Cited by 671



- 2008: Ours much more scalable, graph-based
- 2008: Lim et al., Static birthmarking
- 2008: Zhou et al., Combined approach
- 2008, Roy et al., NiCad
- 2009: Ours for similar code in binaries
- 2009: Juergens et al., ConQat
- 2009: Nguyen et al., ClemanX
- 2010: Lee et al., instant / intelligent code search
- 2010: Hummel et al., index-based
- 2010: Duala-Ekoko et al., clone tracking
- 2011: Tairas et al., sub-clones / clones in DSLs
- 2011: Göde et al., clone changes
- 2011: Higo et al., PDG clones
- 2011: Kim et al., MeCC
- 2012: Dang et al., XIAO clones
- 2012: Rahman et al., code clone smells
- 2012: Hanna et al, Juxtapp for Android apps
- 2012: Nguyen et al., JSync
- 2013: Khoo et al., Rendezvous for binaries
- 2013: Ng et al., Expose: binary clones
- 2013: Crussell et al., Andarwin for Android clones
- 2014: Hu et al., Duet for Android libraries
- 2014: Chen et al., Android clones
- 2014: Lin et al., clonepedia
- 2015: Wang et al., WuKong for Android apps
- 2015: Tian et al., birthmarks on dynamic seqs
- 2015: Wong et al., clocom: clones for comments
- 2016: Sajnani et al., SourcererCC
- 2016: White et al., deep learning for clones
- 2016: Su et al., code relatives
- 2017: Li et al., LibD for Android
- 2017: Cheng et al., cross-languages clones in diffs
- 2017: Palomba et al., code smells
- 2018: Kim et al., FaCoY
- 2018: Ragkhitwetsagul et al., image-based
- 2018: Wang et al., CCAAligner
- 2018: Saini et al., Oreo
- 2018: Zhao et al., DeepSim



Bug Detection

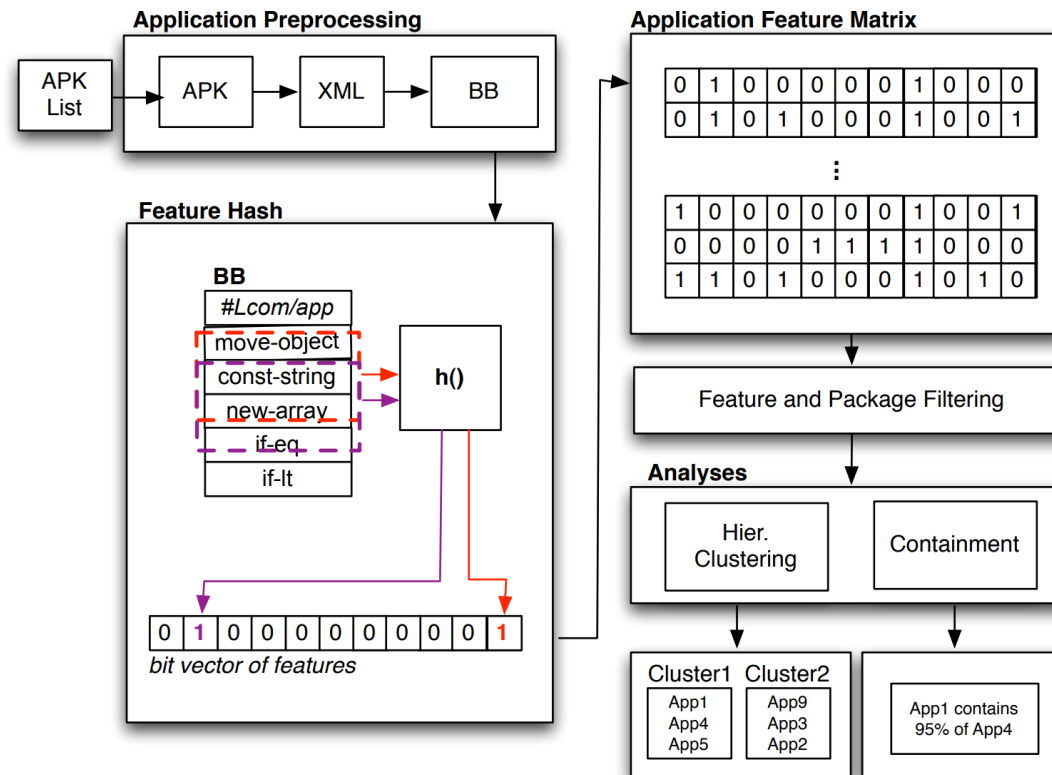
CVE patch

```
...  
+ strcpy_s(dst,src);  
- strcpy(dst,src);  
...
```



```
int boo(char* src)  
{  
    (...)  
    char* dst[10];  
    [...] strcpy(dst,src) [...]  
    (...)  
    return 0;  
}
```

Malware Detection



Refactoring

```
-  if (fContainer == null) {
+  IActionBars actionBar = fContainer.getActionBars();
+  if (actionBars==null&&offset!=0&&!fContainerProvided){
      return Utilities.findActionBars(fComposite, offset);
  }
-  return fContainer.getActionBars();
+  return actionBar;
}
public ISLocator getServiceLocator (int offset2) {
-  if (fContainer == null) {
+  ISLocator sLocator = fContainer.getServiceLocator();
+  if(sLocator == null&&offset2!=0&&!fContainerProvided){
      return Utilities.findSite(fComposite, offset2);
  }
-  return fContainer.getServiceLocator();
+  return sLocator;
```

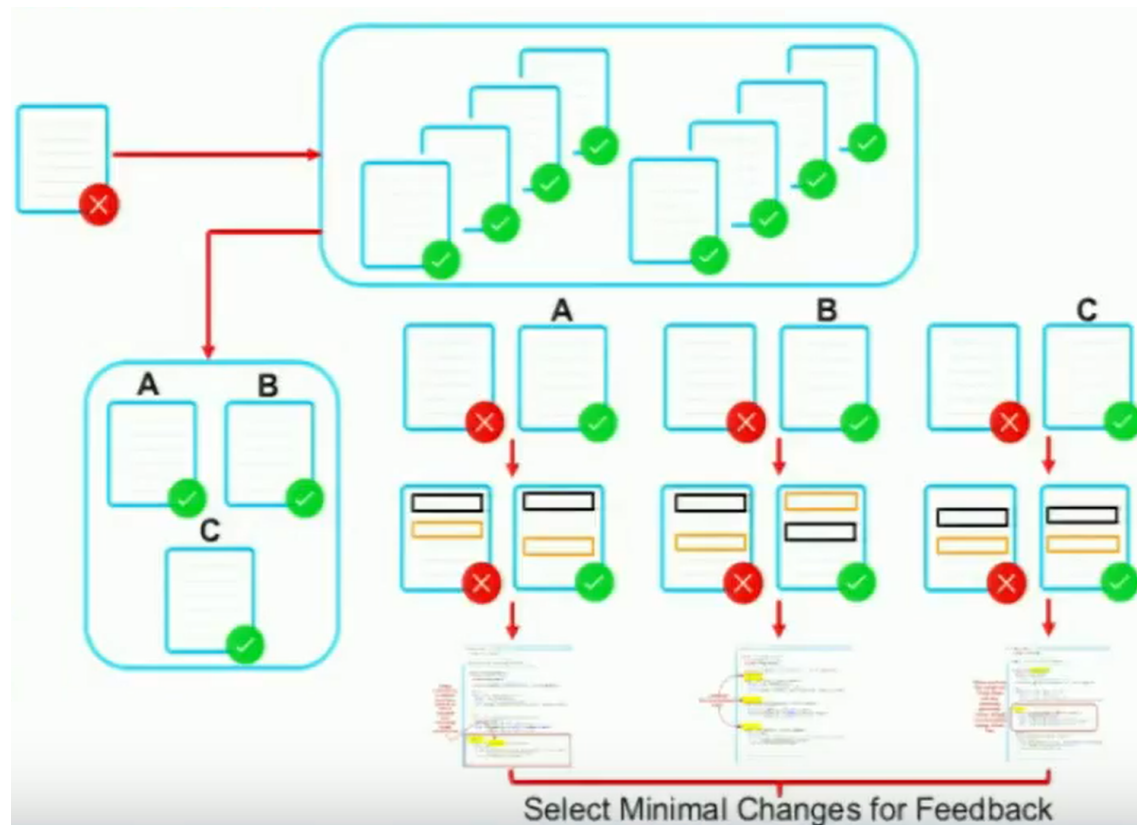
Refactoring

```
- if (fContainer == null) {
+ IActionBars actionBar = fContainer.getActionBars();
+ if (actionBars==null&&offset!=0&&!fContainerProvided){
    return Utilities.findActionBars(fComposite, offset);
}
- return fContainer.getActionBars();
+ return actionBar;
}
public ISLocator getServiceLocator (int offset2) {
- if (fContainer == null) {
+ ISLocator sLocator = fContainer.getServiceLocator();
+ if(sLocator == null&&offset2!=0&&!fContainerProvided){
    return Utilities.findSite(fComposite, offset2);
}
- return fContainer.getServiceLocator();
+ return sLocator;
```



```
T$0 v$0 = fContainer.m$0();
if (v$0==null && v$1!=0 && !fContainerProvided) {
    return Utilities.m$1(fComposite, v$1);
}
return v$0;
```


Feedback Generation



Reflections



High-dimension vector similarity

Eureka moment
during a chat
with a friend



High-dimension vector similarity

LSH was
not enough

High-dimension vector similarity

Hack
internals

LSH was
not enough

High-dimension vector similarity

Hack
internals

LSH was
not enough

Size-based
grouping

High-dimension vector similarity

Hack
internals

LSH was
not enough

Size-based
grouping

Combine with
"normal" hashing

Tree Representation

"Standard"
GCC front-end

Tree Representation

"Standard"
GCC front-end

GCC
Gimple

Tree Representation

"Standard"
GCC front-end

GCC
Gimple

Custom-built
parser generator

False Positives

Probabilistic
Hashing

False Positives

Probabilistic
Hashing

q-level
patterns

False Positives

Probabilistic
Hashing

q-level
patterns

Relevance

False Negatives

Probabilistic
Hashing

False Negatives

Scope

Probabilistic
Hashing

False Negatives

Probabilistic
Hashing

Scope

Granularity

False Negatives

Probabilistic
Hashing

Scope

Semantic
awareness

Granularity

Looking forward

Is Deckard still relevant?

Code Representation

Diff
changes

Symbolic
automata

Gated
graphs

Complex
network

Notions of Similarity & Difference

Operational
semantics

Input-Output

Specification

Description

Beyond Code

Code

Documentation

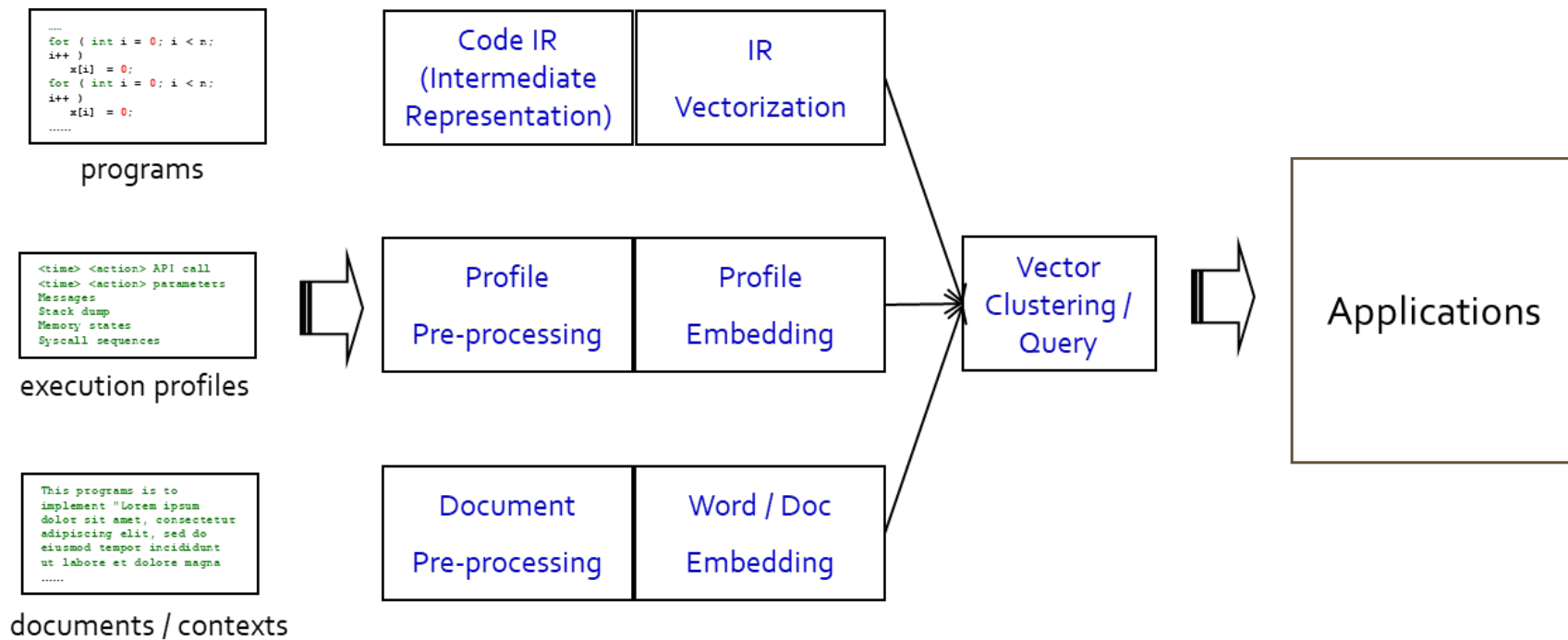
Q&A

Code

**Runtime
Profiles**

Context

Generalization



Looking much further

NATO Software Engineering Conference 1968



```
1 #include <stdio.h>
2 #include <string.h>
3
4 int main(){
5     char *word = "everest";
6     char reverseword[strlen(word)+1];
7     unsigned int letters_remaining = strlen(word);
8     char *wordpointer = &word[strlen(word)-1];
9     int i = 0;
10    while(letters_remaining > 0){
11        reverseword[i++] = *wordpointer--;
12        letters_remaining--;
13    }
14    reverseword[strlen(word)] = '\0';
15    printf("So the reversed word is %s\n",reverseword);
16    return 0;
17 }
```

~

```

1  #ifndef __ASM_ALPHA_FPU_H
2  #define __ASM_ALPHA_FPU_H
3
4  #include <asm/special_insns.h>
5  #include <uapi/asm/fpu.h>
6
7  /* The following two functions don't need trapb/excb instructions
8     around the mf_fpcr/mt_fpcr instructions because (a) the kernel
9     never generates arithmetic faults and (b) call_pal instructions
10    are implied trap barriers. */
11
12    static inline unsigned long
13    rdfpcr(void)
14    {
15        unsigned long tmp, ret;
16
17        #if defined(CONFIG_ALPHA_EV6) || defined(CONFIG_ALPHA_EV67)
18            __asm__ __volatile__ (
19                "ftoit $f0,%0\n\t"
20                "mf_fpcr $f0\n\t"
21                "ftoit $f0,%1\n\t"
22                "itofl %0,$f0"
23                : "=r"(tmp), "=r"(ret));
24        #else
25            __asm__ __volatile__ (
26                "stt $f0,%0\n\t"
27                "mf_fpcr $f0\n\t"
28                "stt $f0,%1\n\t"
29                "ldt $f0,%0"
30                : "=m"(tmp), "=m"(ret));
31        #endif
32
33        return ret;
34    }
35
36    static inline void
37    wrfpcr(unsigned long val)
38    {
39        unsigned long tmp;
40
41        #if defined(CONFIG_ALPHA_EV6) || defined(CONFIG_ALPHA_EV67)
42            __asm__ __volatile__ (
43                "ftoit $f0,%0\n\t"
44                "itofl %1,$f0\n\t"

```

.../linux/arch/alpha/include/asm/fpu.h [Ins] 

(1, 0) [Top/1.8k]

Eval: START

WorkspaceProjectEditAssistantRunGitProfile

EXECdev-machine: clean install#2401:12

Projects Explorer

web-java-spring-petclinic [spring-petclinic]
src
main
java
org.springframework.samples
model
BaseEntity.java
NamedEntity.java
Owner.java
Person.java
Pet.java
PetType.java
Specialty.java
Vet.java

Commands

PetOwnerPetTypeVet

```
86  
87 public Owner getOwner() {  
88     return this.owner;  
89 }  
90  
91 protected void setVisitsInternal(Set<Visit> visits) {  
92     this.visits = visits;  
93 }  
94  
95 protected Set<Visit> getVisitsInternal() {  
96     if (this.visits == null) {  
97         this.visits = new HashSet<Visit>();  
98     }  
99     return this.visits;  
100 }  
101  
102 public List<Visit> getVisits()  
103     List<Visit> visits = getVisitsInternal();  
104     return visits;  
105 }  
97:18
```

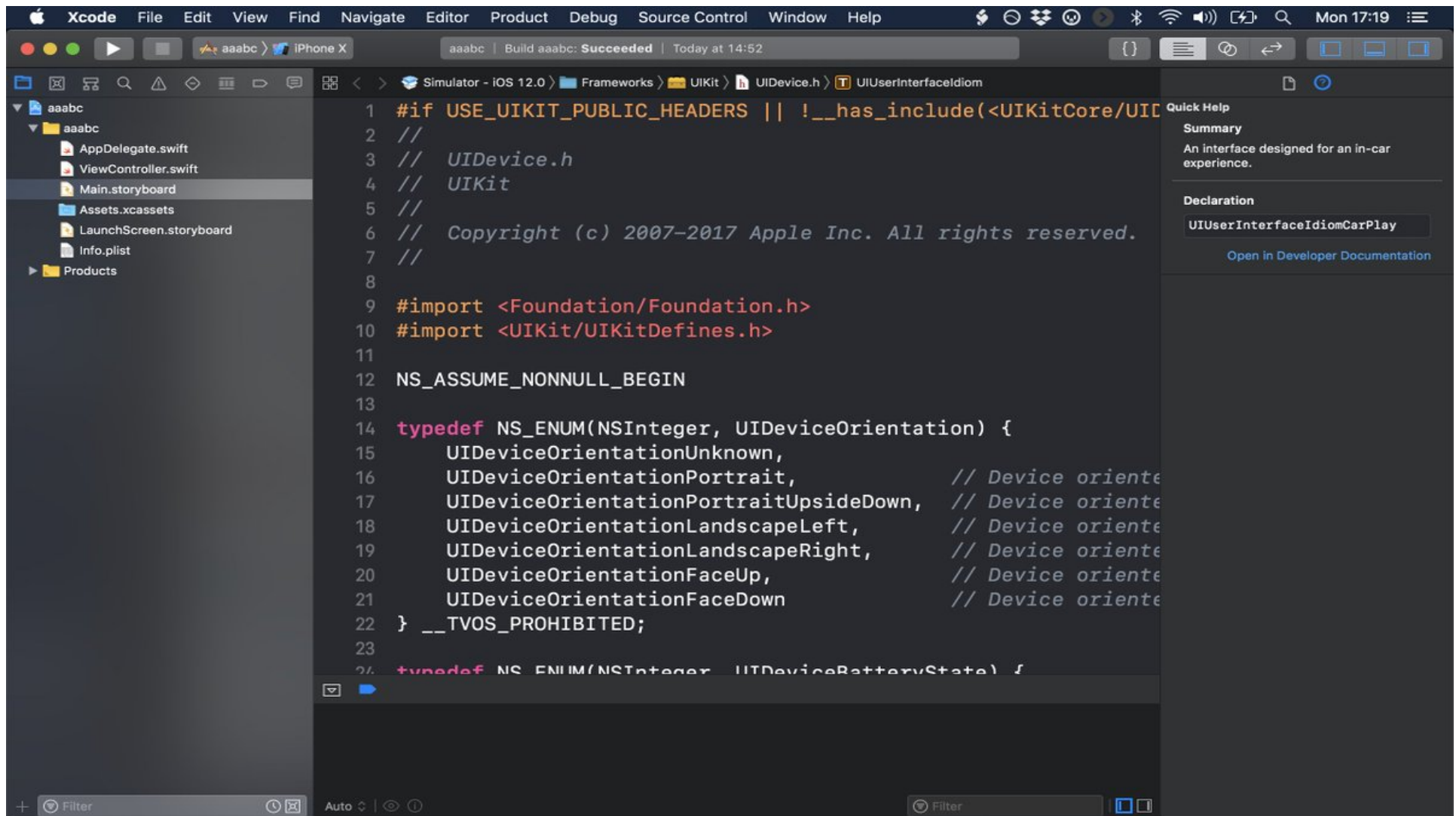
Proposals:
birthDate : DateTime - Pet
id : Integer - BaseEntity
owner : Owner - Pet
type : PetType - Pet
visits : Set<org.springframework.samples.petclinic.model.Visit> - Pet
addVisit(Visit visit) : void - Pet
clone() : Object - Object
equals(Object arg0) : boolean - Object
finalize() : void - Object
getBirthDate() : DateTime - Pet

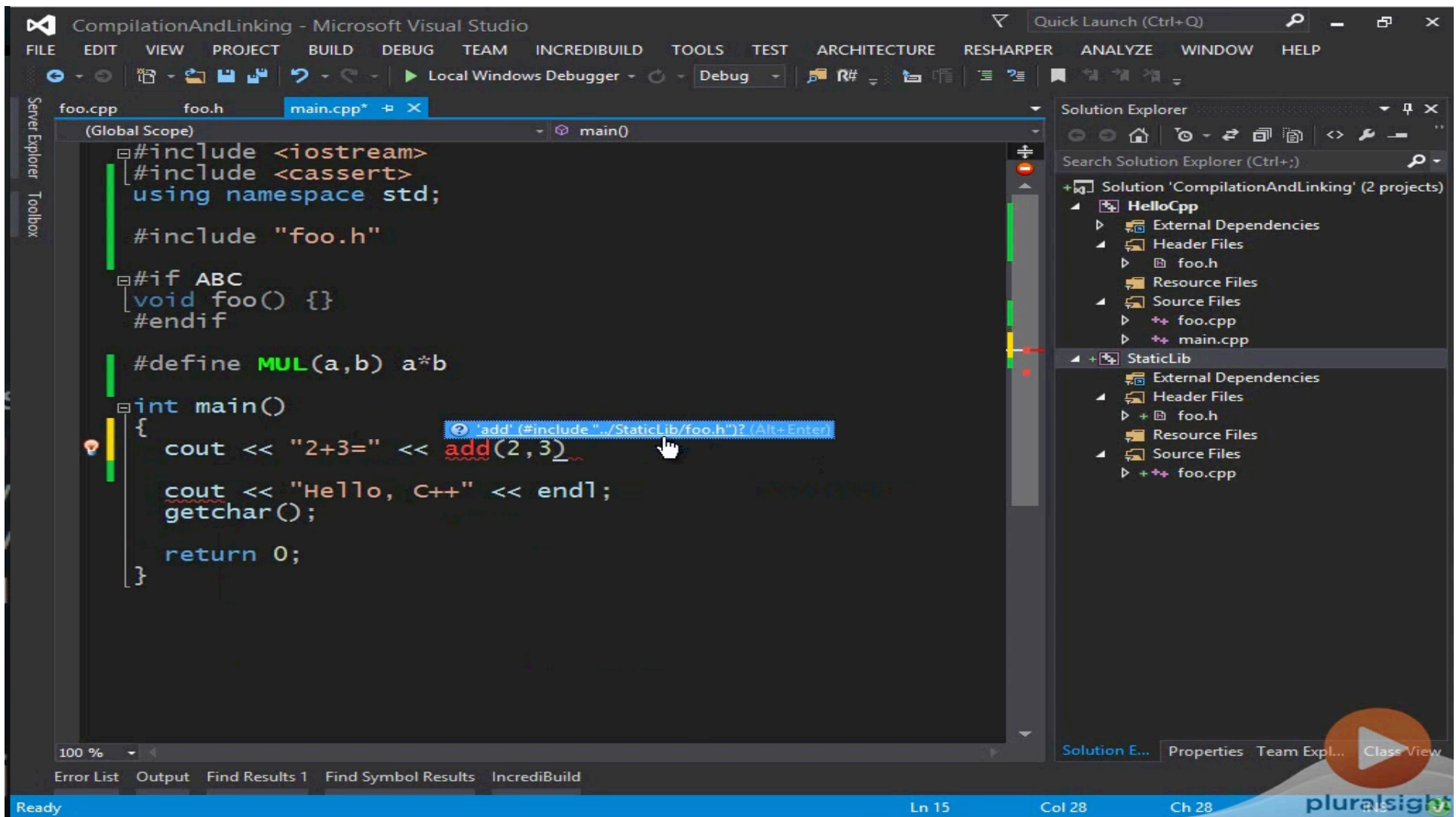
Processes

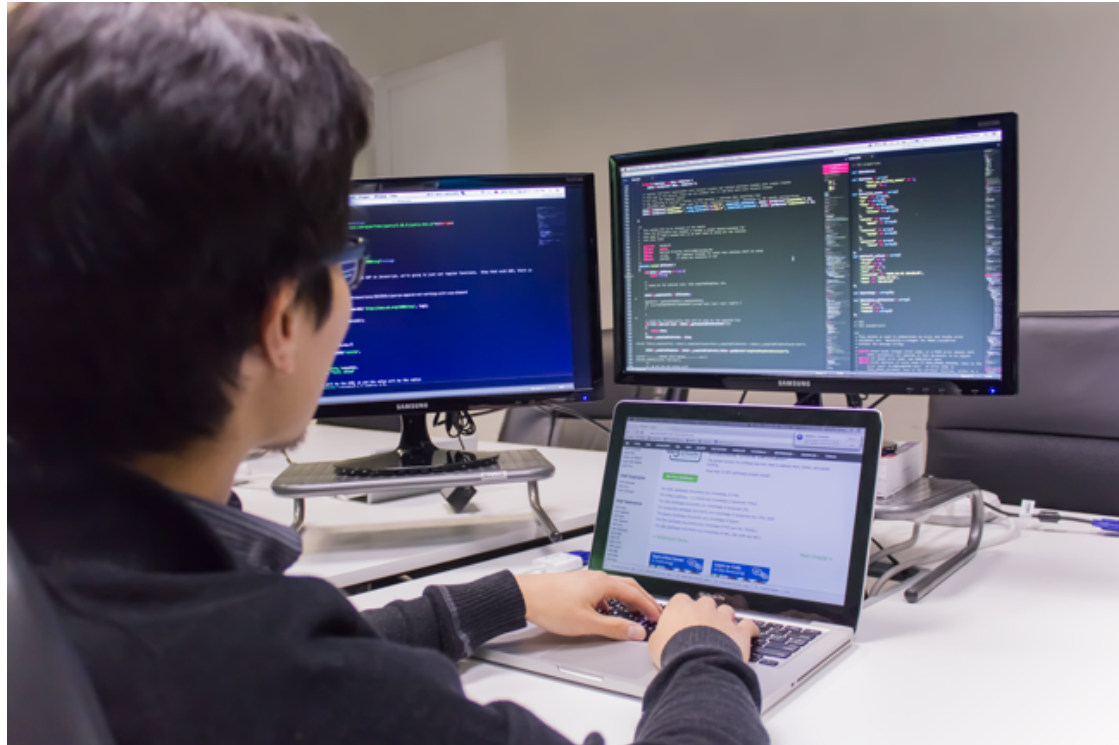
DEVws-machineSSHws-machineTerminalbuild and run

Terminal	drwxr-xr-x	2	root	root	6	Jan	19	16:33	run
build and run	drwxr-xr-x	2	root	root	6	Jan	19	16:33	run
	dr-xr-xr-x	317	root	root	0	Mar	8	08:51	proc
	drwxr-xr-x	3	user	root	4096	Mar	6	18:07	projects
	drwx-----	2	root	root	37	Jan	19	16:33	root
	drwxr-xr-x	6	root	root	127	Mar	8	08:52	run

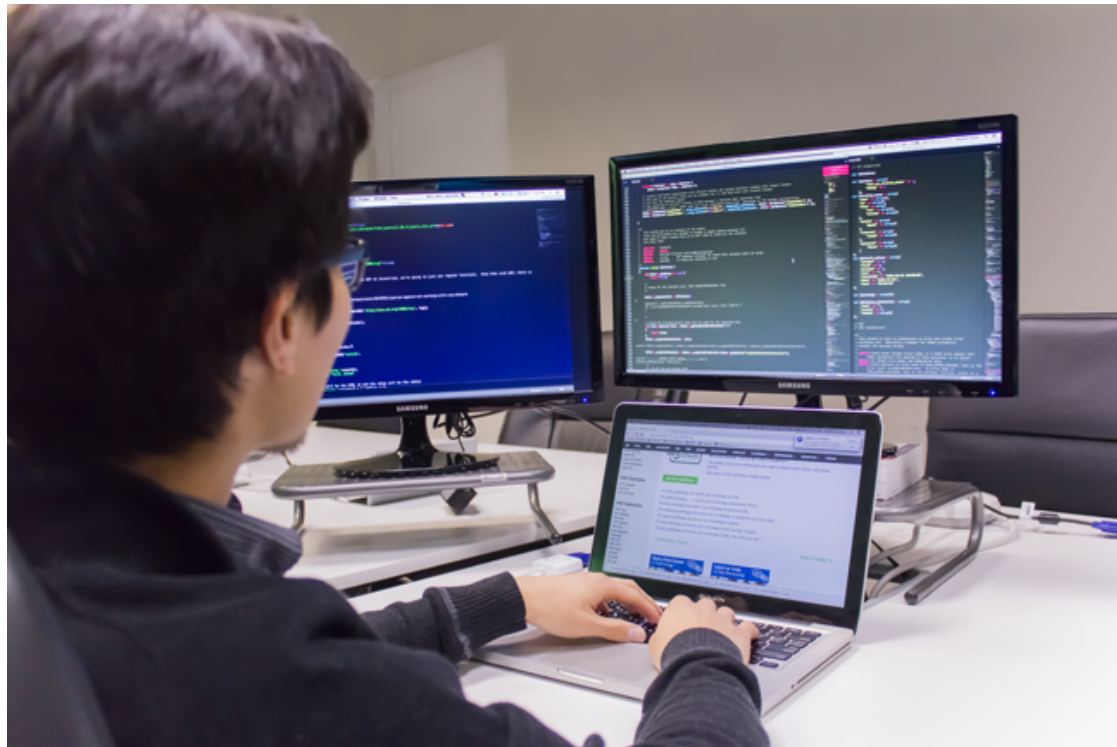
EventsProcesses



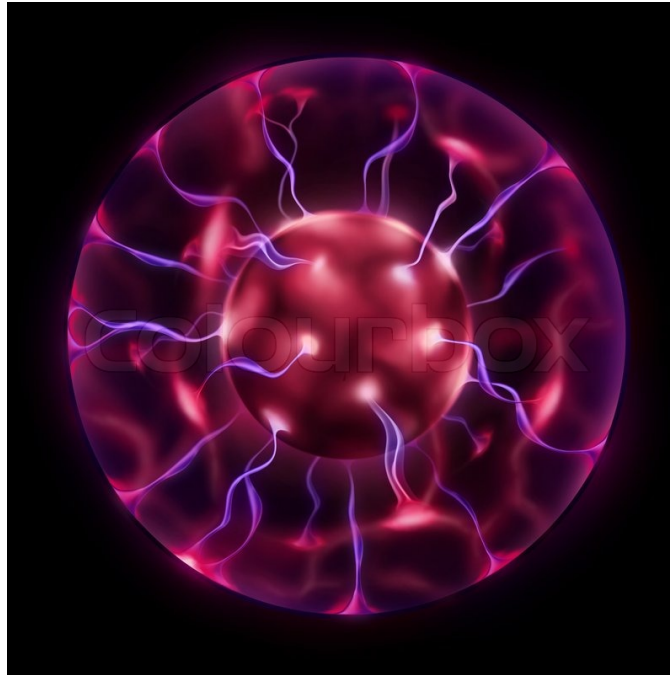




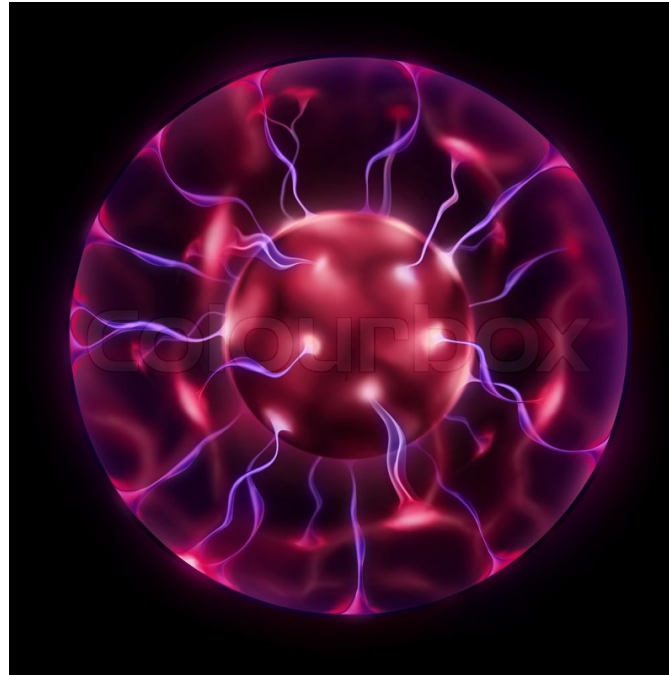
Can we move beyond “coding”?



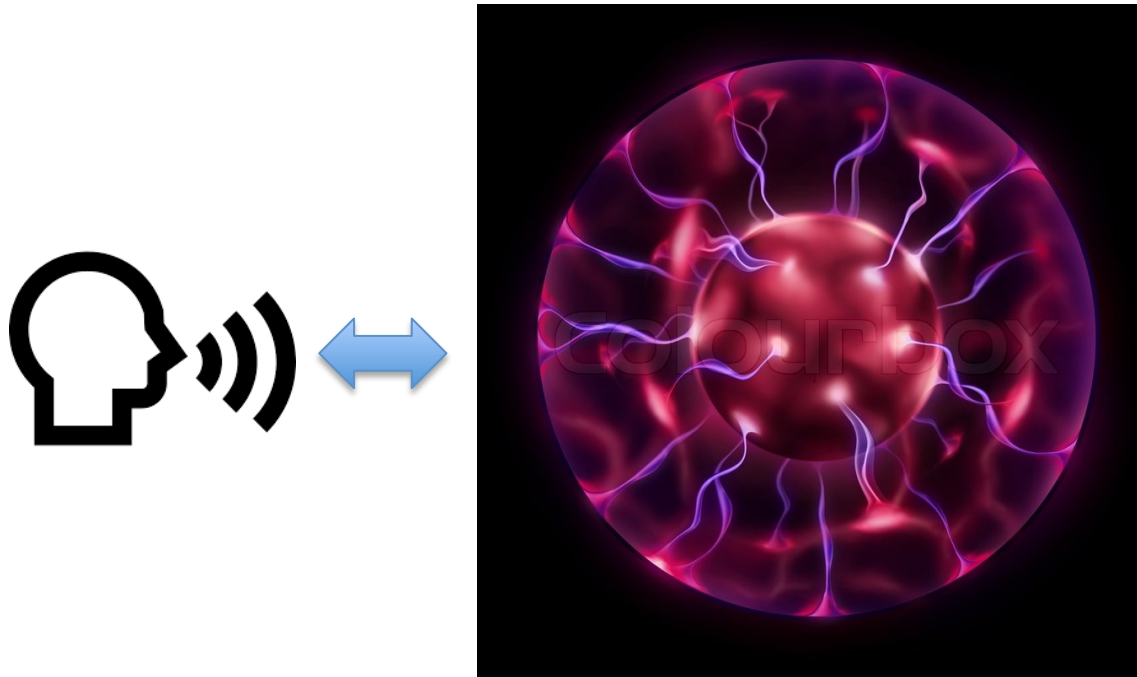
your wish?



Similarity is key



Communicating the wish the hardest





```
function drawTree () {
  var blossomPoints = [];

  resetRandom();
  drawBranches(0, -Math.PI/2, canvasWidth/2, canvasHeight, 30,

  resetRandom();
  drawBlossoms(blossomPoints);
}

function drawBranches (i,angle,x,y,width,blossomPoints) {
  ctx.save();

  var length = tween(i, 1, 62, 12, 3) * random(0.7, 1.3);
  if (i == 0) { length = 107; }

  ctx.translate(x,y);
  ctx.rotate(angle);
  ctx.fillStyle = "#000";
  ctx.fillRect(0, -width/2, length, width);

  ctx.restore();

  var tipX = x + (length - width/2) * Math.cos(angle);
  var tipY = y + (length - width/2) * Math.sin(angle);

  if (i > 4) {
    blossomPoints.push([x,y,tipX,tipY]);
  }

  if (i < 6) {
    drawBranches(i + 1, angle + random(-0.15, -0.05) * Math.PI);
    drawBranches(i + 1, angle + random( 0.15, 0.05) * Math.PI);
  }
  else if (i < 12) {
    drawBranches(i + 1, angle + random( 0.25, -0.05) * Math.PI);
  }
}
```



```

var length = tween(i, 1, 62, 12, 3) + random(0.7, 1.3);
if (i == 0) { length = 107; }

ctx.translate(x,y);
ctx.rotate(angle);
ctx.fillRect(153, -width/2, length, width);

ctx.restore();

var tipX = x + (length - width/2) * Math.cos(angle);
var tipY = y + (length - width/2) * Math.sin(angle);

if (i > 4) {
  blossomPoints.push([x,y,tipX,tipY]);
}

if (i < 6) {
  drawBranches(i + 1, angle + random(-0.15, -0.05) * Math.PI);
  drawBranches(i + 1, angle + random( 0.15,  0.05) * Math.PI);
}
else if (i < 12) {
  drawBranches(i + 1, angle + random( 0.25, -0.05) * Math.PI);
}
}

function drawBlossoms (blossomPoints) {
  var colors = ["#f5ceea", "#e8d9e4", "#f7c9f3", "#ebb4cc", "#f5ceea"];
  ctx.globalAlpha = 0.60;

  for (var i = 0; i < blossomPoints.length; i++) {
    var p = blossomPoints[i];
    for (var j = 0; j < 16; j++) {
      var x = lerp(p[0], p[2], random(0,1)) + random(-10,10);
      var y = lerp(p[1], p[3], random(0,1)) + random(-10,10);

      ctx.fillStyle = colors[Math.floor(random(0,colors.length))];
      ctx.fillCircle(x, y, random(2,5));
    }
  }
}

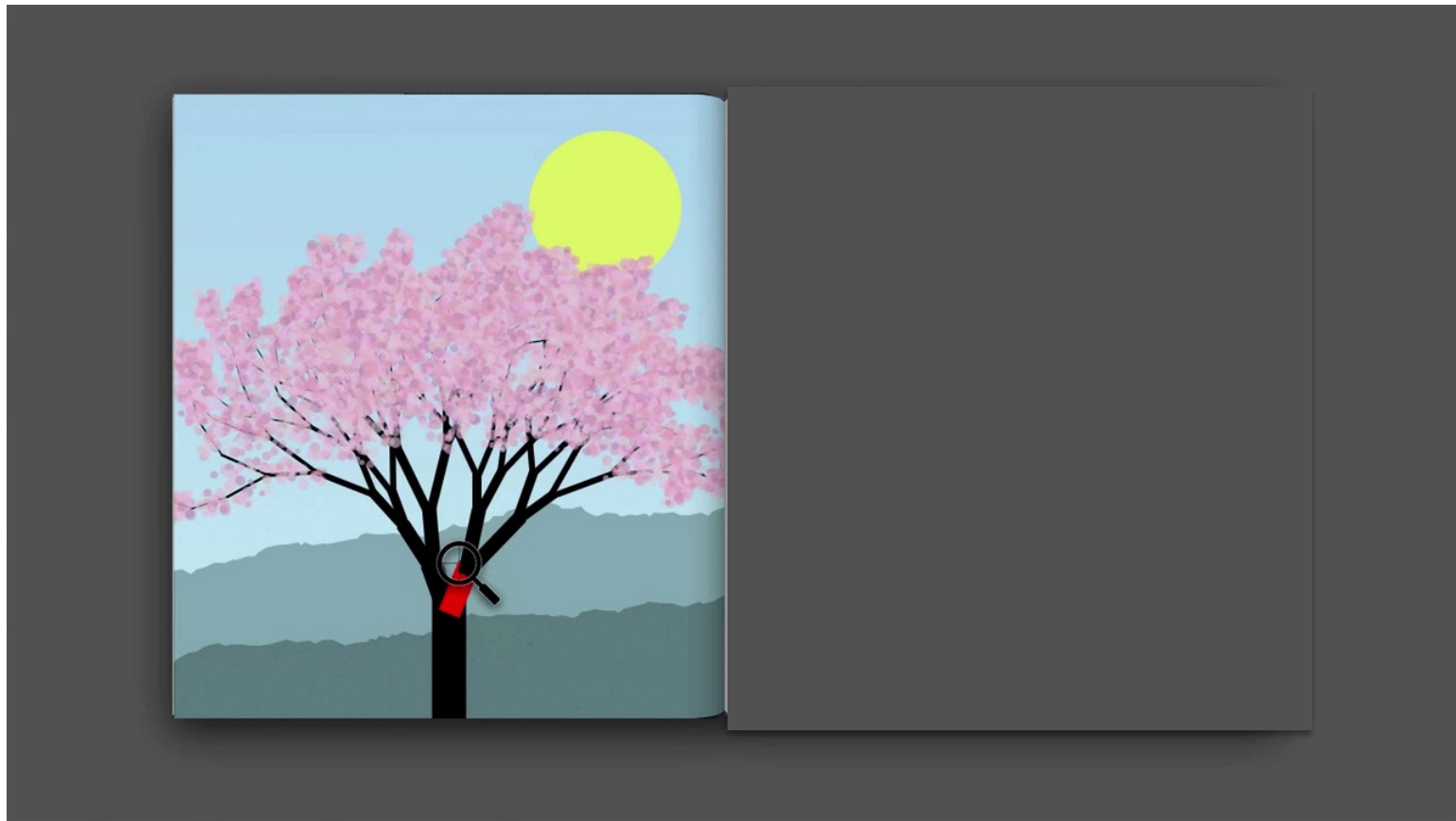
```


Bret Victor: Inventing on Principle





```
function drawTree () {  
  var blossomPoints = [];  
  
  resetRandom();  
  drawBranches(0, -Math.PI/2, canvasWidth/2, canvasHeight, 30,  
  
  resetRandom();  
  drawBlossoms(blossomPoints);  
}  
  
function drawBranches (i,angle,x,y,width,blossomPoints) {  
  ctx.save();  
  
  var length = tween(i, 1, 62, 12, 3) * random(0.7, 1.3);  
  if (i == 0) { length = 107; }  
  
  ctx.translate(x,y);  
  ctx.rotate(angle);  
  ctx.fillStyle = "#000";  
  ctx.fillRect(0, -width/2, length, width);  
  
  ctx.restore();  
  
  var tipX = x + (length - width/2) * Math.cos(angle);  
  var tipY = y + (length - width/2) * Math.sin(angle);  
  
  if (i > 4) {  
    blossomPoints.push([x,y,tipX,tipY]);  
  }  
  
  if (i < 6) {  
    drawBranches(i + 1, angle + random(-0.15, -0.05) * Math.PI);  
    drawBranches(i + 1, angle + random( 0.15,  0.05) * Math.PI);  
  }  
  else if (i < 12) {  
    drawBranches(i + 1, angle + random( 0.25, -0.05) * Math.PI);  
  }  
}
```



Goal is the **object**, not the code

Can we directly
manipulate & explore the object
to express the “wish”?

Can we directly
manipulate & explore the object
to express the “wish”?

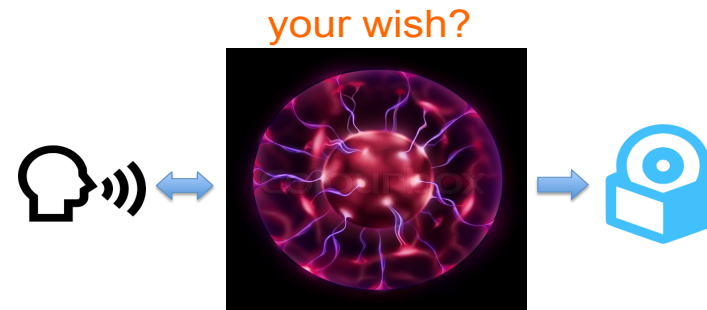
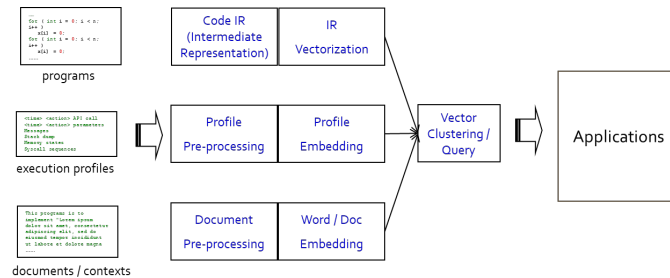
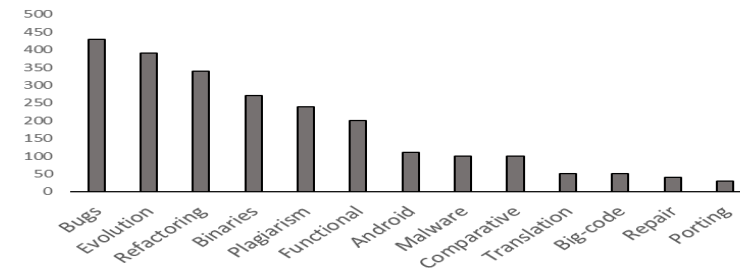
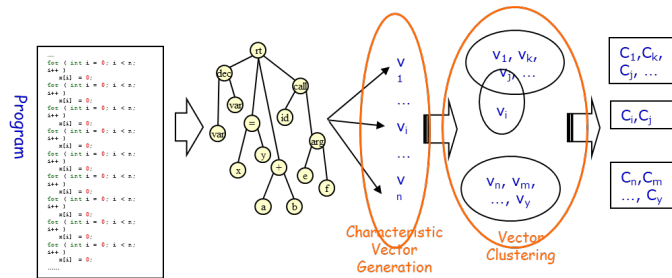
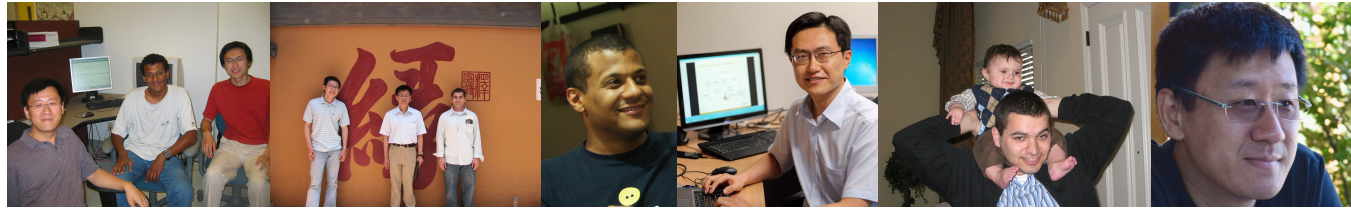
Perhaps via visualization & virtual reality?

We had fun!

We worked hard!

We are grateful!

Software is about similarity!



Thank you!